



## A Probabilistic Analysis Framework for Malicious Insider Threats

Chen, Taolue ; Kammuller, Florian; Nemli, Ibrahim ; Probst, Christian W.

*Published in:*

Proceedings of the third International Conference on Human Aspects of Information Security, Privacy, and Trust (HAS 2015)

*Link to article, DOI:*

[10.1007/978-3-319-20376-8\\_16](https://doi.org/10.1007/978-3-319-20376-8_16)

*Publication date:*

2015

*Document Version*

Peer reviewed version

[Link back to DTU Orbit](#)

*Citation (APA):*

Chen, T., Kammuller, F., Nemli, I., & Probst, C. W. (2015). A Probabilistic Analysis Framework for Malicious Insider Threats. In T. Tryfonas, & I. Askoxylakis (Eds.), *Proceedings of the third International Conference on Human Aspects of Information Security, Privacy, and Trust (HAS 2015)* (pp. 178-189). Springer. Lecture Notes in Computer Science Vol. 9190 [https://doi.org/10.1007/978-3-319-20376-8\\_16](https://doi.org/10.1007/978-3-319-20376-8_16)

---

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

# A Probabilistic Analysis Framework for Malicious Insider Threats<sup>\*</sup>

Taolue Chen<sup>1</sup>, Florian Kammüller<sup>1</sup>, Ibrahim Nemli<sup>2</sup>, and Christian W. Probst<sup>2</sup>

<sup>1</sup> Middlesex University London, {t.chen, F.KammueLLer}@mdx.ac.uk

<sup>2</sup> Technical University Denmark, ibrahimnemli@msn.com, cwpr@dtu.dk

**Abstract.** Malicious insider threats are difficult to detect and to mitigate. Many approaches for explaining behaviour exist, but there is little work to relate them to formal approaches to insider threat detection. In this work we present a general formal framework to perform analysis for malicious insider threats, based on probabilistic modelling, verification, and synthesis techniques. The framework first identifies insiders' intention to perform an inside attack, using Bayesian networks, and in a second phase computes the probability of success for an inside attack by this actor, using probabilistic model checking.

## 1 Introduction

Cyber security considers attacks on organisations from cyber space [5]. While many organisations are well protected against technical attacks, combinations of technical attacks with human factors can be devastating. This integration of human factors and security is important, and extends security to organisational issues and society. This combination has almost replaced the classical “security sciences”, since it is now apparent that in almost all aspects of security the human factor is crucial. However, it is an open challenge how to integrate human behaviour into the design and verification of (secure) systems. A common problem for security analysts is to detect attacks by insiders. Here, more than anywhere, human behaviour needs to be taken into account when designing security systems and monitoring information systems.

In this paper, we present a framework that leverages probabilistic modelling and verification techniques for the analysis of insider threats. As we have shown in previous work [2], insider threat analysis requires the combination of a macro-level view and a micro-level view akin to sociological techniques. This is needed in order to integrate human factors into the context of an infrastructure, like the physical environment of a company and its IT network. We use Bayesian networks to probabilistically model the human disposition for the micro-level analysis to estimate when an actor becomes an insider. Additionally, on the

---

<sup>\*</sup> Part of the research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement no. 318003 (TREsPASS). T. Chen is partially supported by an oversea grant from the State Key Laboratory of Novel Software Technology, Nanjing University.

macro-level we use Markov Decision Processes (MDPs) to model actions of an insider within an organisation’s infrastructure (physical and logical). This two-fold framework provides a tool for the security analyst to *quantitatively* estimate the actual risk of insider threats by an employee of a company at a given moment.

The micro and macro level are represented in our framework as an intentional analysis and a behavioural analysis. The *intentional analysis* (see Section 5) analyses the degree of the intention, in terms of probability, for an employee to be an insider attacker. Once an employee *intends* to be an insider attacker, the *behavioural analysis* (see Section 6) identifies the probability of success, using Probabilistic Model Checking (PMC) to support our analysis. To the best of our knowledge, this is the first quantitative framework to provide a comprehensive analysis of malicious insider threats.

Before presenting this framework, we discuss related work followed by an introduction of the basic concepts of insider threats and probabilistic modelling techniques in Section 3.

## 2 Related work

We base our work on existing taxonomies of insiders [10,3]. In previous work [2], we used Higher Order Logic (HOL) to model insider threats accommodating the view of the insider’s disposition based on these taxonomies, and insider patterns based on real case studies [3]. This logical modelling of insider patterns revealed that HOL allows modelling the human factor with its psychological disposition, the company’s infrastructure including policies, and use theorem proving to prove that certain behaviours lead to policy violations, i.e., insider attacks. However, the need of a company’s security services to quantitatively estimate the risk of an insider attack needs a more detailed analysis like the probabilistic one we present here. The Insider threat patterns provided by CERT [3] use System Dynamics models, which can express dependencies but do not support probabilities quantifying these dependencies nor any of the probabilistic analysis that we propose here. Axelrad et al. [1] have used Bayesian networks for modelling insider threats, and we are currently investigating how their work relates to our first phase. In earlier work [9], we have used EXASyM [11] to model and analyse attacks based on the infrastructure of a company expressed as a graph in the acKlaim calculus and using the PRISM model checker. There, attacks are based on the probabilities of actors’ movements in the infrastructure corresponding to a random walk. Here, we embed this previous analysis tool to provide the tooling for the second part of our framework, but significantly extended to Markov Decision Processes (as opposed to merely Markov Chains) as used in [9].

## 3 Preliminaries

The behavioural and psychological aspects of actors are related to their personal profile in many ways. These can be seen as the antecedents or key initial factors to understanding an individual’s propensity to perform an attack. Nurse *et*

al. [10] have identified eight elements that may be especially useful in modelling and analysing this aspect of insider threats. These are the precipitating events (catalyst), an individual's general personality characteristics, historical behaviour, concrete psychological state in a situation, attitudes towards work, skill set, opportunity, and lastly, motivation to attack. All of these values are hard to measure, but if present, can be used to compute probabilities of the occurrence and success of actions.

We represent these probabilities in *Bayesian networks* (BNs) [6], which correspond to the micro level view of our framework. To construct these in general, one has to collect many possible observations that may be relevant to the problem and determine what subset of those observations is worthwhile to model, and then organise the observations into variables having mutually exclusive and collectively exhaustive states [4].

BNs are a graphical model that encodes a probabilistic relationship among variables of interests. In general, a BN for a set of variables  $\mathcal{X} = \{X_1, \dots, X_n\}$  is a tuple  $(S, P)$  where  $S$  is a *directed acyclic graph* (DAG) that encodes the set of conditional independence assertions on variables in  $\mathcal{X}$ , and  $P$  is a set of local probability distributions associated with each variable. These two components together define a joint probability distribution for  $\mathcal{X}$ . The nodes in  $S$  are in one-to-one correspondence with the variables in  $\mathcal{X}$ . We usually use  $X_i$  to denote both the variable and its corresponding node, and  $Pa(X)$  to denote the set of parents of the node  $X$  in  $S$  as well as the variables corresponding to those parents. In  $S$ , the absence of edges between two nodes (variables) encodes conditional independencies (of the two variables). In particular, given the structure  $S$ , the joint probability distribution for  $\mathcal{X}$  is given by

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i \mid pa_i)$$

where  $pa_i$  denotes the parents of node  $X_i$  in  $S$ , and  $P$  are the distributions corresponding to the term  $p(\cdot \mid \cdot)$ .

On the macro level, we use Markov Decision Processes (MDPs) to identify potentially successful insider threats:

**Definition 1 (MDP).** A Markov Decision Process  $\mathcal{M} = (S, s_0, A, \mathbf{P}, AP, L)$ , where

- $S$  is a set of states with  $s_0 \in S$  being the initial state;
- $A$  is a set of actions;
- $\mathbf{P} : S \times A \times S \rightarrow [0, 1]$  is the transition probability function such that for all states  $s \in S$  and actions  $a \in A$ ,  $\sum_{t \in S} \mathbf{P}(s, a, t) \in \{0, 1\}$ .
- $AP$  is a set of atomic propositions;
- $L : S \rightarrow 2^{AP}$  is the labelling function.

For any state  $s \in S$  and action  $a$ , if  $\sum_{t \in S} \mathbf{P}(s, a, t) = 1$ , then we say the action  $a$  is enabled in  $s$ .

## 4 Framework

The overall aim of our analysis is to estimate the probability that an employee of an organisation (conceived as the insider) launches a successful insider attack. We emphasise that, in this work, we only address *intentional, malicious* insider threats, meaning that the insider consciously acts as an attacker. In contrast, the analysis of *accidental* insider threats – where the insider might be social engineered by a malicious outsider – is not addressed and is left as future work. Overall, our framework consists of two components:

- The intentional analysis provides a quantitative measure for the risk that a particular employee may reach the tipping point and turn into a malicious insider; and
- The behavioural analysis estimates where an insider could successfully launch an attack in a company’s infrastructure. This is influenced by the personal characteristics of the attacker, for example, the attacker’s skill to break a lock or succeed in social engineering the secretary.

Our framework is quantitative in terms of probabilities, which depend on various factors, typically including a *personal profile* and a *type of insider threats*. A *personal profile* comprises a number of factors influencing the person’s behaviour; the profile does not describe the behaviour, it is a prediction of likely behaviour, based for example on tests, profiles, and observed behaviour:

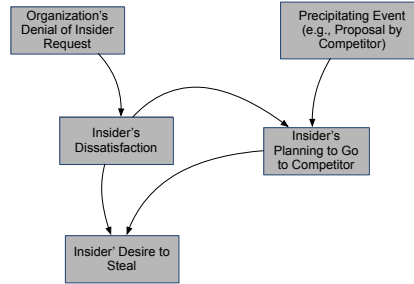
- Individuals’ personality characteristics. The personnel department will often be able to provide an estimation via some “characteristics test” commonly used in psychology;
- Psychological state;
- Attitude towards work: should be easily evaluated, and formalised as a variable in the interval  $[0, 1]$ ; or
- Skill set: for instance whether one can break a locker, whether one has good knowledge of CCTV, etc.

The personal profile will be used for both the intentional analysis and the behavioural analysis below. We also note that the personal profile is *time-dependent* and should be updated regularly as some attributes might become invalidated, jeopardising the precision of the analysis.

The CERT Guide [3] identifies three main types of insider threats, fraud, theft of intellectual property, and sabotage. Evidently, these types influence both steps of our analysis; for instance, fraud and theft require different skills, and so the probability of success differs.

## 5 Intentional Analysis

In this section, we provide the first component of our framework, *i.e.*, the intentional analysis of potential insider threats. The main aim is to have an estimation on the *degree* that an employee of an organisation intends to launch an insider



**Fig. 1.** Example for System Dynamics model for Entitled Independent [3].

attack. Formally, this is expressed as the *probability* that an employee decides to become an insider attacker. At this stage we do *not* address whether finally the attack is accomplished or not; this is the main subject of the behavioural analysis.

The main part of the analysis is based on Bayesian networks. Modelling using BNs is widely considered as an art, and requires sufficient domain-specific knowledge. However, for insider threats there exist collections of so called *patterns* [3], which turn out to facilitate our modelling significantly. In particular, the *System Dynamics* modelling method is exploited. In this methodology abstract variables define a taxonomy of insider threat cases. Graphically, these variables are presented in square boxes. A *solid* arrow from a box containing variable  $a$  to one containing variable  $b$  indicates that an increase of  $a$  implies an increase of  $b$ . A *dashed* arrow represents the inverse relationship, *i.e.*, an increase of  $a$  implies a *decrease* of  $b$ . An example pattern for “entitled independent” is given in Fig. 1.

The system dynamics, which only gives a high-level, qualitative description, still yields a useful starting point of the intentional analysis. Essentially, we substantially extend the system dynamics by introducing a *quantitative* description, which is considerably more precise and useful. Note that the causalities, quantification, or qualification are encoded as a Bayesian network. In light of this, a crucial step of our methodology is to translate the system dynamics to BN.

The BN has the same graph structure as the system dynamics model, which we assume to be acyclic.<sup>3</sup> In general, for each node in the system dynamics, we introduce a random variable. We usually have the following cases:

1. For *events* that might happen or not, the corresponding random variable is governed by the *Bernoulli distribution*, *i.e.*, a random variable which takes value 1 with success probability  $p$  and value 0 with failure probability  $q = 1 - p$ ; A typical case of this kind is the precipitating event in Fig. 1;

<sup>3</sup> We note that in practice, occasionally the cycles exist, which we abstract to their strongly connected component as a single node, thus obtaining a proper DAG. We might be able to use Markov logic network to directly encode a system dynamic with cycles, but this is left as the future work for simplicity.

2. For quantities over a finite domain, for instance, the psychological state which might take values from *happy*, *depressed*, *disgruntled*, *angry*, *stressed*, as well as the type of motivation which might take values from *financial*, *political revenge*, *fun*, *competitive-advantage*, *power*, *peer-recognition*, we usually introduce discrete random variables with corresponding outcomes as the domain of the quantity under consideration;
3. For quantities of continuous nature, for instance, dissatisfaction in Fig. 1, in principle, we can introduce a continuous random variable, say the *degree* of dissatisfaction, as a value in  $[0, 1]$ . In practice, we usually apply discretisation to  $[0, 1]$  to have a partition of  $[0, 1]$ .

As the next step, we must specify the conditional probabilities among the introduced random variables. It is worth noting that the concrete probabilities are difficult to obtain; however, this is not the main concern of the current paper which solely aims to establish the basic framework.

We emphasise that such a model should be *parameterised*, for example with the type of threat. The reason is that for different threat types, the introduced random variables should vary, and probably more importantly, the conditional probabilities differ.

Once the BN is established, the next step is to analyse it. In general, we abstract the degree of intention as a value  $\text{Int} \in [0, 1]$ , and the analysis computes the probability that the degree of intention falling into interval  $I$  exceeds  $\theta$ , i.e.,  $\Pr[\text{Int} \in I] \geq \theta$ . This is a typical task of *prediction*. Another kind of analysis is *explanation*, for instance, the analyst might be interested in knowing, once the degree of intention  $\text{Int} \geq \theta$ , what is the most likely cause?

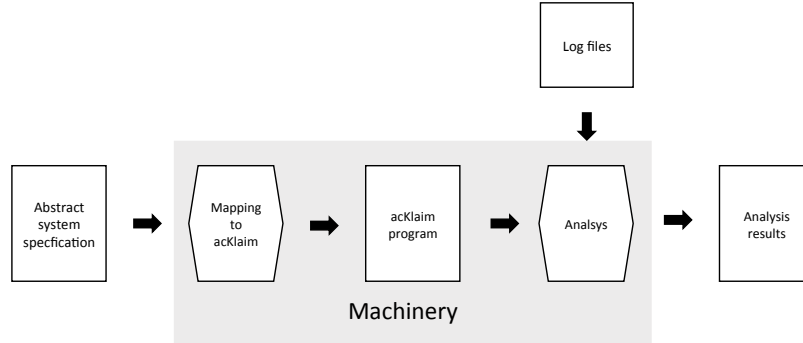
## 5.1 Example

To illustrate the intentional analysis framework, let's consider the case study of *Entitled Independent* depicted in Fig. 1.

The events “organisation denial of insider request” and “precipitating event” are of type (1), and thus are governed by the Bernoulli distribution. For instance we have that  $\Pr(\text{precipitating event} = 1) = 0.1$  meaning that with probability 0.1, the employer receives a job offer from a different company. The “insider’s dissatisfaction” is of type (3), and thus we consider the *degree* of dissatisfaction  $\beta \in [0, 1]$  and introduce a probabilistic density function (pdf)  $f_\beta$  to specify the distribution of  $\beta$ . However for computation efficiency, we usually prefer to stay in the discrete model, so we could partition  $[0, 1]$  into  $[0, 0.1)$ ,  $[0.1, 0.2)$ ,  $\dots$  and specify the probabilities  $p_1, \dots, p_{10}$ . The intuition is that, say, “the probability that the dissatisfaction degree being from 0.5 to 0.6 is  $p_5$ ”.

Alternatively, we can define “low”, “mediate”, “high” for the dissatisfaction degree, which could correspond to  $[0, 0.3)$ ,  $[0.3, 0.7)$ ,  $[0.7, 1]$  respectively. This is also the case for “Insider’s desire to steal”, and “Insider’s planning to go to competitor”, which practically takes value from “yes”, “no”, “not sure”.

As the next step, we need to specify the conditional probability. For the most interesting case regarding an insider’s desire to steal, this could simplistically be



**Fig. 2.** Workflow of the behavioural analysis framework [9]

formulated as a table, where X denotes “Insider’s dissatisfaction”, and Y denotes “Insider’s planning to go to competitor”.

X	Y	steal (H)	steal (M)	steal (L)
High	Yes	0.6	0.3	0.1
High	Not sure	0.2	0.3	0.5
High	No	0.0	0.1	0.9
⋮	⋮	⋮	⋮	⋮

## 6 Behavioural analysis

In this section, we provide the second component of our framework, i.e., the behavioural analysis of insider threats. Given the infrastructure of the organisation and a personal profile, this analysis estimates the probability of successful insider attacks. Here the infrastructure of the organisation may refer to the physical locations relevant to the insider threats, their access control policies, etc.

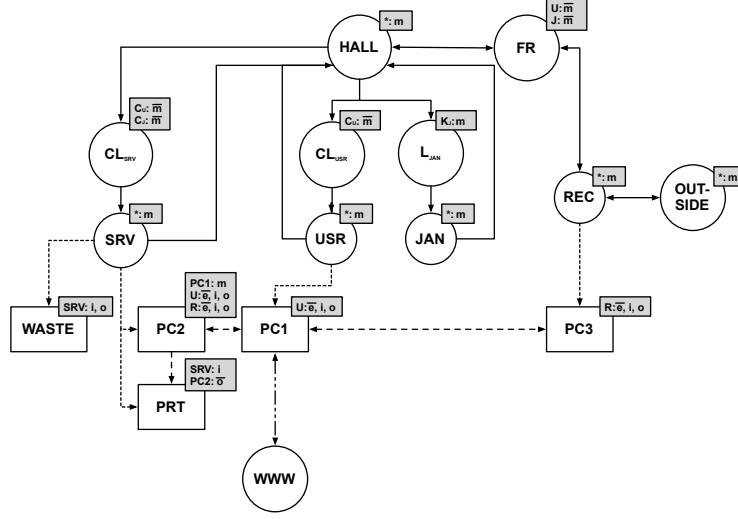
Fig. 2 illustrates the workflow, which consists of the following four main steps:

- Step 1: Model the infrastructure in the *abstract system specification*;
- Step 2: Map the abstract system specification to the *acKlaim process calculus* and generate the *transition system*;
- Step 3: Translate the *transition system* into *Markov decision processes* by annotating the transitions with probabilities; and
- Step 4: Perform behavioural analysis by verification of the Markov decision process.

We now elaborate these steps in details.

**Step 1.** The abstract system model is specified as a collection of mathematical constructs that is used to create an abstraction of a real-world system. This





**Fig. 3.** A simple example system and its representation as a graph including actors, networks, locations, and policies in acKlaim [11]

abstraction makes it possible to model physical localities, interconnected computers, actors that can move around in the physical localities, and data that can be carried by actors or left at both computers or localities. On top of this there is a fine-grained access control mechanism that limits the mobility of actors, and protects sensitive data. Figure 3 shows an example of a physical model and its representation in acKlaim [11].

It should be noted that in the abstract system, there is no means for modelling dynamic behaviours of actors, but only means of specifying what the initial structure of the system looks like – a static representation of the system model. The dynamic behaviour of the model is supported when the model is mapped to acKlaim in which the semantics of acKlaim uses the abstract system to evaluate the effects of actors movement.

We note further that the location might be *physical* or *logical*. Physical locations such as HALL in Fig. 3 are self-explained. In contrast, logical locations provide a valuable means to model human aspects of insider threat. For instance, it is useful to capture the “secretary” by-pass. Indeed, it is well recognised that a typical scenario of insider threats is that the attacker obtains privilege (e.g., entering certain restricted areas, obtaining the master key) by having special personal relationship with secretary-like persons such as personal assistant of the CEO, the receptionist.

**Step 2.** We use the acKlaim process calculus to model an organisation and its actors as a graph of locations and actors. In order to explore insider behaviour in organisational models, we use an abstract view on policy formalisations and analysis: policies describe prerequisites for actions to be granted to actors given by

pairs of predicates (conditions) and sets of enabled actions. We integrate policies into the infrastructure, providing an organisational model where policies reside at locations and actors are adorned with additional predicates to specify their credentials. In Figure 3, the policies are given in grey boxes. For example, the policy  $\{U: \text{elog}, i, o; \text{pc1}: e;\}$  attached to the (virtual) node `pc2` describes that the user `U` can evaluate processes on `pc2` but his actions are being logged, `U` can input from and output to `pc2`, and `pc1` can evaluate on `pc2`.

To support automated analysis, we have developed a *system specification language* [9] to specify acKlaim system models in text files, which are the input for the analyses. The acKlaim models are translated automatically into a transition system [9, Section 3.5].

**Step 3** Once the transition system is generated, we augment it to obtain a Markov decision process.

**Definition 2 (TS).** *A transition system is a tuple  $(S, s_0, \rightarrow)$  where  $S$  is a set of states with  $s_0 \in S$  being the initial state, and  $\rightarrow \subseteq S \times S$  is a transition relation.*

We note that each *state* of the TS denotes a location in Step 1. We identify a subset of *terminal* states  $F \subseteq S$ . Intuitively, these terminal states denote the places where the insider attack is actually happening. In the framework, we consider the following probabilities:

- For each state (being physical or logical), there is an *entering probability* specifying the probability for an actor (insider) to access that location by any means. In practice, this depends on various different factors, including the access control policy the organisation adopted and the personal profile. Formally, the entering probability is defined as  $p_e : S \setminus F \rightarrow [0, 1]$  where  $F$  is the set of terminal states. Evidently,  $p_e$  must satisfy some constraints. In the simplest case, if the actor is allowed to enter by the access control policy, the entering probability  $p_e = 1$ . However, even if access is not granted by the access privilege, there is still a certain probability to enter (e.g., by breaking the lock of the door which depends on the skill set of the person under consideration).
- For each terminal state, there is a *successful probability* specifying the probability that an actor manages to accomplish the attack after entering the terminal location, formally defined by  $p_s : F \rightarrow [0, 1]$ .
- For each state, we also consider the probability of being caught, i.e., when the insider attempts to perform the attack. We consider the four combinations (1) successful attack, undetected; (2) successful attack, detected; (3) failed attack, undetected; (4) and failed attack, detected. Formally, we define two functions  $p_c : S \rightarrow [0, 1]$  specifying the probability of being detected in each state. Intuitively, these probabilities depend on, for instance, the presence of surveillance.

By assuming that the event of being successful and the event of being caught are independent (which is a reasonable assumption in practice), one can derive the probabilities for the aforementioned combinations of events.

As before, the model is parameterised. These parameters are used to model various factors which would impact the probabilities  $p_e(\cdot)$ ,  $p_s(\cdot)$ , and  $p_c(\cdot)$  introduced above substantially. A typical case is the daytime vs night mode; breaking in during night time might be easier, so the entering probability must be higher. However, at night, some server might be shut down, so the success probability might be lower if one wants to download a confidential file from the internal server. Moreover, the model should be parameterized with personal profiles to account for attacker skills.

With these probabilities at hand, given the TS  $(S, s_0, \rightarrow)$  obtained from the previous step, we can define an MDP  $\mathcal{M} = (S', s'_0, A, \mathbf{P}, AP, L)$  as follows:

- The state space of the MDP  $\mathcal{M}$ ,  $S' = (S \cup \{\text{succ}, \text{fail}\}) \times \{\checkmark, \times\}$ ; and the initial state  $s'_0 = (s_0, \checkmark)$ .
- $A = \{e_{s,t} \mid (s, t) \subseteq \rightarrow\} \cup \{\text{commit}\}$ ;
- For each state of the form  $(s, \star)$  where  $s$  is *not* a terminal state and  $\star \in \{\checkmark, \times\}$ , we introduce an action  $e_{s,t}$  for each edge  $(s, t)$  in the TS. We define the resulting probability distribution  $\mathbf{P}(s, e_{s,t}, \cdot)$ , written  $\mu_{s,t}$  by:

$$\begin{aligned}\mu_{s,t}(t, \checkmark) &= p_e(t) \cdot (1 - p_c(t)) \\ \mu_{s,t}(t, \times) &= p_e(t) \cdot p_c(t) \\ \mu_{s,t}(\text{fail}, \checkmark) &= (1 - p_e(t)) \cdot (1 - p_c(t)) \\ \mu_{s,t}(\text{fail}, \times) &= (1 - p_e(t)) \cdot p_c(t)\end{aligned}$$

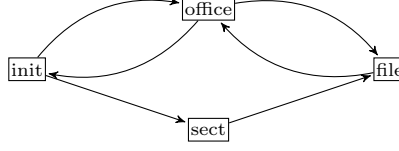
- For each terminal state  $s$ , there is only one action **commit** enabled at  $s$ , and we define the resulting distribution  $\mathbf{P}(s, \text{commit}, \cdot)$ , written  $\mu_{\text{commit}}$ , by:

$$\begin{aligned}\mu_{\text{commit}}(\text{succ}, \checkmark) &= p_s(s) \cdot (1 - p_c(t)) \\ \mu_{\text{commit}}(\text{succ}, \times) &= p_s(s) \cdot p_c(t) \\ \mu_{\text{commit}}(\text{fail}, \checkmark) &= 1 - p_s(s) \cdot (1 - p_c(t)) \\ \mu_{\text{commit}}(\text{fail}, \times) &= 1 - p_s(s) \cdot p_c(t)\end{aligned}$$

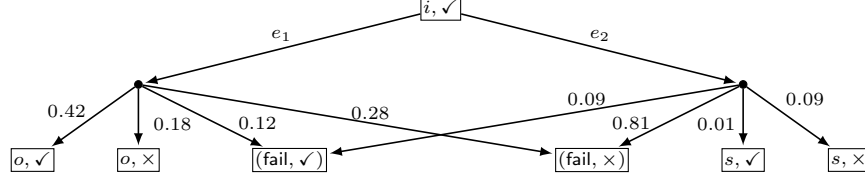
- $(\text{succ}, \star)$  and  $(\text{fail}, \star)$  for  $\star \in \{\checkmark, \times\}$  are *absorbing* states whose transitions do not affect the analysis. We hence omit the definition here.

Intuitively, the  $\checkmark$  means that the insider is *not* caught, while the  $\times$  denotes that the insider is caught; **succ** denotes success of the insider, while **fail** denotes failure. We remark that the definitions of atomic propositions and labelling functions depend on the properties one wants to analyse. We postpone their definitions to the next step.

*Example 1.* We give a (simplified) example to illustrate the construction of MDPs. The transition system is depicted in Fig. 4. From the init state, the insider could go to the office and from there to the file state, or go to the sect state (bribes the secretary) and from there to the file state as well. For clarity, we only depict the MDP corresponding to the state  $(\text{init}, \checkmark)$ . For this state, there are two enabled actions  $e_{\text{init}, \text{office}}$  (denoted by  $e_1$  in Fig. 5) and  $e_{\text{init}, \text{sect}}$  (denoted by  $e_2$  in Fig. 5).



**Fig. 4.** Transition system.



**Fig. 5.** The fragment of MDP for state init;  $i$ ,  $o$ ,  $s$  abbreviate init, office, sect respectively.

**Step 5.** For the last step, we analyse the obtained MDPs by standard probabilistic model checking techniques. The number of interest is the maximum probability that the insider steals a confidential file without being caught. Note that the insider has different strategies, for instance, selecting where to go from a physical location, or trying to social engineer other actors. The insider’s goal is to maximise the success probability, whereas from the organisation’s perspective, a worst-case scenario should be considered.

Such a problem boils down to the problem of computing the maximum probability to reach the state  $(succ, \checkmark)$ . Formally, one can introduce  $AP = \{Succ, NCaught\}$  and a labelling function  $L(succ, \checkmark) = \{Succ, NCaught\}$ . The logical formula

$$\mathbf{P}^{\max=?}[\Diamond(Succ \wedge NCaught)]$$

and the probabilistic model checker PRISM is able to return the maximum probability, as well as the corresponding strategy of the insider to achieve this probability. By such an analysis, the organisation can identify the potential weakness of the infrastructure, and carry out necessary security improvement.

## 7 Conclusion

Insider threats are hard to capture in a systematic way. Extending on our earlier work on representing insiders and behaviour with Higher Order Logic, we have outlined a framework for identifying malicious insider threats in system models using probabilistic model checking. Using System Dynamics, this approach captures the behaviour of insiders, and models both their intent or risk of turning malicious as well as the risk of an insider action succeeding.

In future work we plan to investigate the relation to attacker profiles, budgets [7,8], and skill sets, threats posed by collaborating insiders, and especially the threat posed by accidental insider threats.

## References

1. E. T. Axelrad, P. J. Sticha, O. Brdiczka, and J. Shen. A bayesian network model for predicting insider threats. In *2013 IEEE Security and Privacy Workshops*, pages 82–89, Los Alamitos, CA, USA, 2013. IEEE Computer Society.
2. J. Boender, M. G. Ivanova, F. Kammüller, and G. Primiero. Modeling human behaviour with higher order logic: Insider threats. In *STAST’14*. IEEE, 2014. co-located with CSF’14 in the Vienna Summer of Logic.
3. D. M. Cappelli, A. P. Moore, and R. F. Trzeciak. *The CERT Guide to Insider Threats: How to Prevent, Detect, and Respond to Information Technology Crimes (Theft, Sabotage, Fraud)*. SEI Series in Software Engineering. Addison-Wesley Professional, 1 edition, Feb. 2012.
4. D. Heckerman. A tutorial on learning with bayesian networks. In *M. Jordan, ed. Learning in Graphical Models*. MIT Press, 1999.
5. R. Kissel. Glossary of key information security terms. Technical Report NISTIR 7298 Revision 2, National Institute of Standards and Technology, 2013.
6. D. Koller and N. Friedman. *Probabilistic Graphical Models - Principles and Techniques*. MIT Press, 2009.
7. A. Lenin and A. Buldas. Limiting adversarial budget in quantitative security assessment. In R. Poovendran and W. Saad, editors, *Decision and Game Theory for Security*, volume 8840 of *Lecture Notes in Computer Science*, pages 155–174. Springer International Publishing, 2014.
8. A. Lenin and A. Buldas. Limiting adversarial budget in quantitative security assessment. In *5th International Conference, GameSec 2014*, pages 155–174, 2014.
9. I. Nemli. Using acclaim and prism to model and analyse insider threats. Master’s thesis, DTU Copenhagen, 2015. Available [http://www2.imm.dtu.dk/pubdb/views/edoc\\_download.php/6864/pdf](http://www2.imm.dtu.dk/pubdb/views/edoc_download.php/6864/pdf).
10. J. R. C. Nurse, O. Buckley, P. A. Legg, M. Goldsmith, S. Creese, G. R. T. Wright, and M. Whitty. Understanding Insider Threat: A Framework for Characterising Attacks. In *WRIT’14*. IEEE, 2014.
11. C. W. Probst and R. R. Hansen. An extensible analysable system model. *Information Security Technical Report*, 13(4):235–246, Nov. 2008.